# CUNY SPS DATA 621 - CTG5 - HW2

*Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Koh*

*March 13, 2019*

## Contents

Deliverables should use R functions and the other packages to generate the `classification metrics` for the provided data set.

**1. Download the classification output data set (attached in Blackboard to the assignment).**

**2. The data set has three key columns we will use:**

- class: the actual class for the observation
- scored.class: the predicted class for the observation (based on a threshold of 0.5)
- scored.probability: the predicted probability of success for the observation

**Use the table() function to get the raw confusion matrix for this scored dataset. Make sure you understand the output. In particular, do the rows represent the actual or predicted class? The columns?**

|   | 0 | 1 |
|---|---|---|
| 0 | 119 | 30 |
| 1 | 5 | 27 |

- rows = predicted, cols = actual
- 1 is positive, 0 is negative

**3.-8.: Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns (3)Acurracy; (4)Error rate; (5)Precision; (6)Sensitivity(recall); (7)Specificity; (8)F1 score of the predictions. Verify that you get an accuracy and an error rate that sums to one.**

**3. Accuracy**

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{True Positives} + \text{True Negatives} + \text{False Positives} + \text{False Negatives}}$$

**4. Error Rate**

$$\text{Error Rate} = \frac{\text{False Positives} + \text{False Negatives}}{\text{True Positives} + \text{True Negatives} + \text{False Positives} + \text{False Negatives}}$$

**5. Precision**

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

**6. Sensitivity (or Recall)**

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

**7. Specificity**

$$\text{Specificity} = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$$

**8. F1 Score**

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}}$$

```
## [1] 1
```

| accuracy | error.rate | precision | sensitivity | specificity | f1 |
|---|---|---|---|---|---|
| 0.8066298 | 0.1933702 | 0.7986577 | 0.9596774 | 0.4736842 | 0.8717949 |

**9. Before we move on, let's consider a question that was asked: What are the bounds on the F1 score? Show that the F1 score will always be between 0 and 1.(Hint: If 0 < a < 1 and 0 < b < 1 then ab < a)**

Given the formula for the F1 score. . .

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}}$$

. . . it can be determined the values of Precision and Sensitivity are of utmost importance.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \ , \text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Considering the extremes can lead to a more whole answer regarding the bounds on the F1 score. First, if we assume TP has a value of 0. . .

$$\text{Precision} = \frac{0}{0 + \text{FP}} \ , \text{Sensitivity} = \frac{0}{0 + \text{FN}}$$

$$\text{Precision} = \frac{0}{\text{FP}} \ , \text{Sensitivity} = \frac{0}{\text{FN}}$$

$$\text{Precision} = 0 \ , \text{Sensitivity} = 0$$

Next, if it is assumed TP has a value of 1. . .

$$\text{Precision} = \frac{1}{1 + \text{FP}} \ , \text{Sensitivity} = \frac{1}{1 + \text{FN}}$$

Depending on the values of FP and FN for Precision and Sensitivity respectively, this will always result in a value greater than 0 and less than or equal to 1. This is the case for all non-zero values of TP.

Since it has been proven that Precision and Sensitivity will have values between 0 and 1 inclusively, the same treatment may be applied to the F1 Score.

If Precision and Sensitivity are both 0. . .

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}}$$

$$\text{F1 Score} = \frac{2 \times 0 \times 0}{0 + 0}$$

$$\text{F1 Score} = \frac{0}{0}$$

$$\text{F1 Score} = 0$$

If Precision and Sensitivity are both 1. . .

$$\text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}}$$

$$\text{F1 Score} = \frac{2 \times 1 \times 1}{1 + 1}$$

$$\text{F1 Score} = \frac{2}{2}$$

$$\text{F1 Score} = 1$$

On each extreme the F1 Score has a value of 0 and a value of 1. That means if Precision and Sensitivity should fall between 0 and 1, the F1 Score should also have a value that falls between 0 and 1.

*In short, the bounds of the F1 score are 0 and 1.*

**10. Write a function that generates an ROC curve from a data set with a true classification column (class in our example) and a probability column (scored.probability in our example). Your function should return a list that includes the plot of the ROC curve and a vector that contains the calculated area under the curve (AUC). Note that I recommend using a sequence of thresholds ranging from 0 to 1 at 0.01 intervals.**

A receiver operating characteristic (ROC) curve plots the true positive rate (sensitivity) against the false positive rate (1 - specificity) for all of the possible cutoff values.

$$\text{TPR} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{FPR} = \frac{\text{False Positive}}{\text{False Positive} + \text{True Negative}}$$

The ROC curve on the figure below was created using build function (Appendix Chapter 10). The Trajectory of the curve between (0.0) and the 60% threshold is steep, indicating that the sensitivity is increasing at a greater rate than decrease in specificity. However, when the sensitivity is greater than 80%, there is a more significant decrease in specificity than the gain in sensitivity.
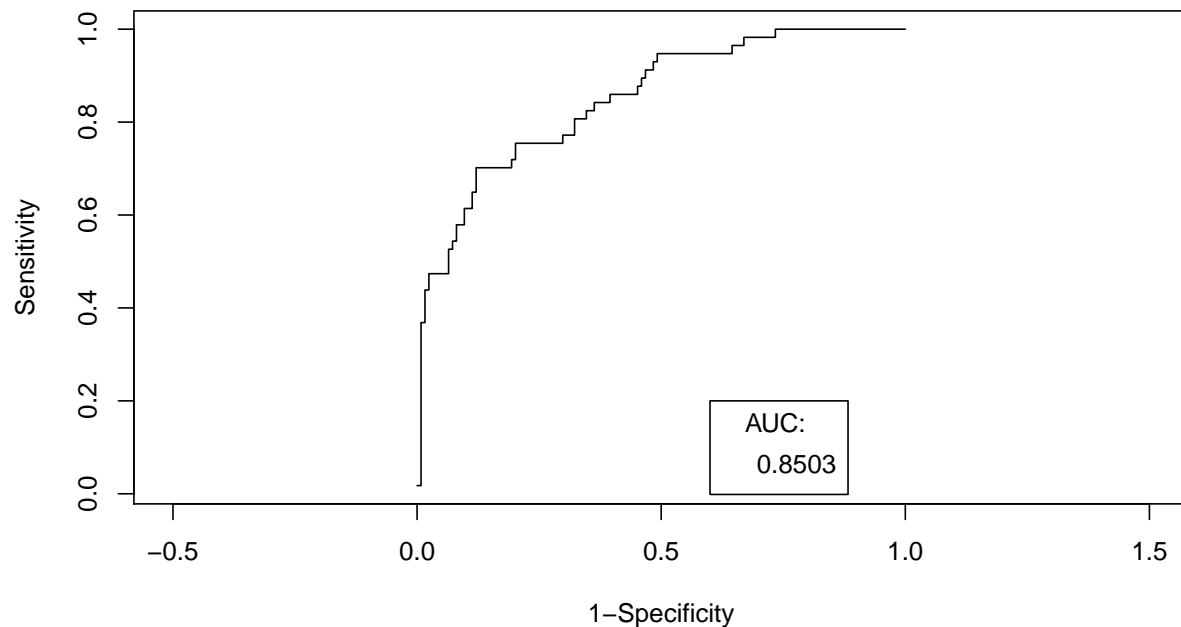
Figure 1: ROC curve using function

The area under the ROC curve (AUC) is a measure, of how well a parameter can be distinguished between two diagnostic groups. Computing the area under the curve is one way to summarize it in a single value. Our model has a good calculated AUC: 0.8503. It's near 1, which means it has good measure of class separation.

**11. Use your created R functions and the provided classification output data set to produce all of the classification metrics discussed above.**

**12. Investigate the caret package. In particular, consider the functions confusionMatrix, sensitivity, and specificity. Apply the functions to the data set. How do the results compare with your own functions?**

The caret ("Classification and REgression Training") package is a toolset for model training. Its function `confusionMatrix` provides a shortcut to cross-tabulate observed and predicated classes with their statistics, just like the custom function created in the preceding questions.

`confusionMatrix` takes as input the actual and predicted class for each observation - both as factors with same set of levels (the approach we used), or as a table. After computing the prevalence of true and false positives and negatives, the function outputs a list of statistics (including several the custom function was not asked to examine, for which formulae are included below):

- Accuracy
- 95% confidence interval and p-value
- No information rate (rate at which correct if just selected majority class) and its p-value
- Kappa coefficient (a measure agreement between actual and predicted classes comparing to probability of chance agreement)
- McNemar's test p-value
- Sensitivity
- Specificity

- PPV (Positive Predicted Value)

$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

- NPV (Negative Predicted Value)

$$\frac{\text{True Negatives}}{\text{True Negatives} + \text{False Negatives}}$$

* Prevalence

$$\frac{\text{True Positives} + \text{False Negatives}}{\text{True Positives} + \text{False Positives} + \text{True Negatives} + \text{False Negatives}}$$

- DetectionRate

$$\frac{\text{True Positives}}{\text{True Positives} + \text{True Negatives} + \text{False Positives} + \text{False Negatives}}$$

- DetectionPrevalence

$$\frac{\text{True Positives} + \text{False Positives}}{\text{True Positives} + \text{True Negatives} + \text{False Positives} + \text{False Negatives}}$$

- Balanced Accuracy

$$\frac{\text{Sensitivity} + \text{Specificity}}{2}$$

The caret packages `ConfusionMatrix` yielded identical results to the statistics computed by the custom function to seven decimal places.

| accuracy | error.rate | precision | sensitivity | specificity | f1 |
|---|---|---|---|---|---|
| 0.8066298 | 0.1933702 | 0.7986577 | 0.9596774 | 0.4736842 | 0.8717949 |

| accuracy | error.rate | precision | sensitivity | specificity | f1 |
|---|---|---|---|---|---|
| 0.8066298 | 0.1933702 | 0.7986577 | 0.9596774 | 0.4736842 | 0.8717949 |

**13. Investigate the pROC package. Use it to generate an ROC curve for the data set. How do the results compare with your own functions?**

The ROC curve on the figure below was created using pRoc Package. Comparing this curve to the one we created using the function, we can see that both curves look identical.

We also got matching result for the area under the curve (AUC): 0.8503.
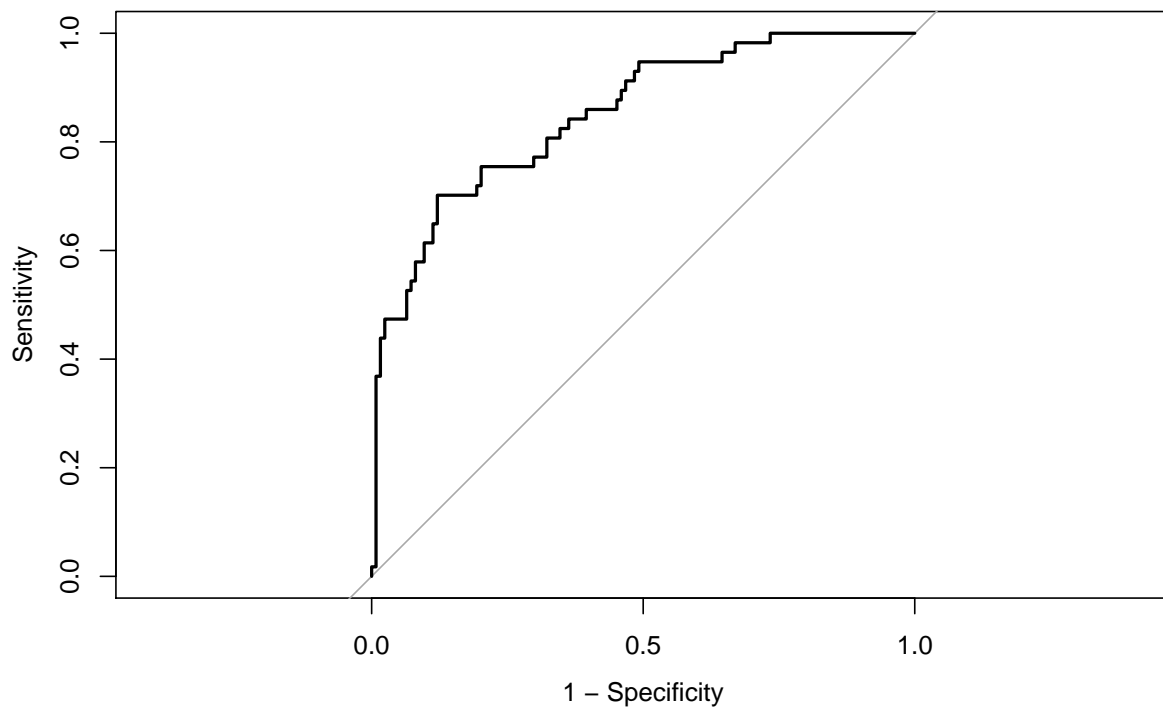
Figure 2: ROC curve pROC package

Using pROC package we can also compute the confidence interval (CI) of a ROC curve. By default, it's 95%. Computed confidence interval for our curve is 0.7905-0.9101.